

Integrating Clouds and Cyberinfrastructure: Research Challenges

Manish Parashar, Rutgers University, parashar@rutgers.edu

Geoffrey Fox, Indiana University, gcf@indiana.edu

Kate Keahey, Argonne National Laboratory

Alan Sill

Scott Brim

Michael Nelson

Aggressive Cloud computing technology development has resulted in many multiple classes of Cloud services that provide attractive solutions for many different types of business applications. It is expected that Cloud services will join more traditional research cyber infrastructure (CI) components, such as high-performance computing system, clusters and Grids in supporting scientific exploration and discovery. It is clear from current research that there are real benefits in using Clouds and Cloud computing abstractions as part of a hybrid cyber infrastructure to support CDS&E, for example, to simplify the deployment of applications and the management of their execution, improve their efficiency, effectiveness and/or productivity, and provide more attractive cost/performance ratios. Furthermore, Clouds and Cloud computing abstractions can support new classes of algorithms and enable new applications formulations, which can potentially revolutionize CDS&E research and education. However, before CDS&E can fully realize the potential benefits of a hybrid cyber infrastructure that integrates Cloud services, several research issues remain. The objective of this report is to explore these research challenges. Note that the discussion below is from the CDS&E perspective and is complementary to more general Cloud research challenges.

The discussion in this report is based on two reports. The first is by Gannon and Fox¹, in which they reviewed and classified applications suitable for Clouds. The second is by Parashar et al.², in which they explored how a hybrid HPC/Grid + Cloud cyber infrastructure can be effectively used to support real-world science and engineering applications, presented illustrative scenarios, and discussed limitations and research challenges. Furthermore, the report is the result of discussions as part of MAGIC meeting in September 2011 and April and May 2012.

1. Algorithms and Application Formulations for Clouds

Most attempts to directly port a conventional HPC application to a Cloud platform fail. The challenge is to think differently and rewrite the application to support the new computational and programming models. It is important to understand what are the key

¹G. Fox, D. Gannon, "Cloud Programming Paradigms for Technical Computing Applications," Technical Report, <http://grids.ucs.indiana.edu/ptliupages/publications/Cloud%20Programming%20Paradigms.pdf>, 2012,

²M. Parashar, M. AbdelBaky and I. Roderio, "Cloud Paradigms and Practices for CDS&E," Technical Report, 2012, <http://cometcloud.org>, 2012.

differences between HPC and Cloud usage modes from an application perspective, and where the available trade-offs are addressed and how.

A key attribute of Clouds is on-demand access to elastic resources, i.e., applications programmatically access more or less resources as they evolve, to meet changing needs. Such a capability can have a significant impact on how algorithms are developed and applications are formulated. For example, the execution of an application no longer has to constrain itself to a fixed set of resources that are available at runtime and can grow or shrink its resource set based on the demands of the science – the science can drive the scale and type of resource involved based on, for example, the levels of refinement required to resolve a solution feature, or the number of ensembles that need to be run to quantify the uncertainty in a solution, or the type of online analytics services that need to be dynamically composed into the application workflow. Understanding how CDS&E applications can effectively utilize Clouds and Cloud abstractions, as part of a hybrid cyberinfrastructure, to enable new practices and levels of scientific insights remains a research challenge.

Research is also needed to explore the meaningful science, engineering and business application scenarios that can take advantage of such hybrid infrastructure, such as data intensive, data-driven, sensor-based, high-throughput, etc. For example, a meaningful HPC plus Cloud use case may consist of simulations with online data analytics/visualization. In such a scenario, exposing the ability to modify goals/configurations based on data analytics feedback to the user will be critical to ensuring impact on the science. For example, in data-intensive computations, the use of feature tracking might allow the scientist to adjust application parameters based on the analysis of meaningful features using a public Cloud, where the analysis can be performed in a timely manner due to shorter resource provisioning times compared to a high-end HPC system. We believe that such meaningful scenarios will exist in all areas of CDS&E.

Attributes of Clouds such as elasticity, on-demand provisioning, multi-tenancy and virtualization can lead to research challenges related to validating correctness of execution, fault-tolerance and fault-tolerant application formulations, reproducibility, and provenance tracking and management. While reproducibility can be improved by encapsulating application context within a virtual machine image, the variability in the execution environment (for example, due to multi-tenancy) can become an issue. Clouds also present new requirements and challenges for debugging and testing.

There are important classes of applications that need special attention and have special research challenges. Examples are:

- Biomedical and bioinformatics applications, where cloud architecture brings special challenges in the area of privacy.
- Sensor-webs, where the elastic nature of Clouds is well suited for the often bursty nature of sensor data.
- Big data applications, for example those based on new MapReduce or Iterative MapReduce environments, result in broad research areas addressing programming and storage challenges. Latter include SQL and NOSQL models and the reconciliation of distributed data and centralized cloud computing.

2. Programming Models, Abstractions and Systems

A key research challenge is developing appropriate programming abstractions and language extensions that can enable CDS&E applications to simply and effectively take advantage of the elastic access to resources and services during application formulation. Furthermore, it may be necessary to define constraints (for example, budgets, data privacy, performance, etc.) to regulate the elasticity, and the programming abstractions may provide support for expressing these constraints so that they can be enforced during execution. Similarly, such annotations can also define possible adaptations, which could then be used to increase performance, manageability and overall robustness of the application. Example annotations include “dynamically increase the assigned resources in order to increase the resolution of a simulation under certain convergence constraints” or “modify convergence goals to avoid failure or guarantee completion time”. The Cloud service model can also lead to interesting services specialized to CDS&E that provide entire applications or applications kernels as a service (i.e., SaaS). Furthermore, and arguably more interestingly, it can also export specialized platforms for science as a services (i.e., PaaS), which encapsulate elasticity and abstract of the underlying hybrid infrastructure.

Programming abstractions are also necessary to enable applications to effectively integrate CI and Clouds and support hybrid models of execution. Similarly, specialized programming support for applications classes can be very effective, such as, for example, the MapReduce or Iterative MapReduce models and system for “Big Data” applications.

Finally, tools for model checking, (probabilistic) verification, configuration management, debugging, coordinated execution, etc., are important components of programming systems for Clouds, given issues such as jitter, performance variance, noisy-neighbor, multi-tenancy, etc.

3. Middleware Stacks and Services, Management Policies, and Economic Models

Middleware services will need to support the new CDS&E applications formulations and services enabled. A key research aspect will be the autonomic management and optimization (multiple objectives including performance, energy, cost, reliability, etc.) of application execution through cross-layer application/infrastructure adaptations. It will be essential for the middleware services to be able to adapt to the application’s behavior as well as system configuration, which can change at run time, using the notion of elasticity at the application and workflow levels. Furthermore, appropriate services are necessary to be able to provision different types of resources on demand. For example, if we target HPC as a Cloud and HPC plus Cloud approaches on the NSF funded cyber-infrastructure such as XSEDE, Open Science Grid and FutureGrid along with commercial Clouds such as Amazon EC2 or Microsoft Azure, autonomic provisioning and scheduling techniques, including Cloud bursting will be necessary to support hybrid usage modes. Finally, monitoring, on-line data analytic for proactive application/resource management and adaptation techniques will be essential as the scale and complexity of both the applications and hybrid infrastructure grows.

There are many research areas related the two sections above. These include:

- Scheduling models optimized for Cloud and hybrid Cloud+CI usage modes such as those discussed above.
- Support for the dynamic (on-demand) federation of Clouds, the linkage of private and public Clouds and of Clouds and CI, and Cloud bursting.
- Optimizing the performance of virtualized systems (compute, communication, network, storage). This includes new reduction primitives, polymorphic implementation on different systems with for example, exploitation of high performance networks as in classic MPI research.
- Autonomic management and optimization (cross-layer, multiple objectives including performance, energy, cost, reliability, etc.). Monitoring, on-line data analytic for proactive application/resource management.
- Interoperability and integration of cloud storage models/solutions. This includes new storage models including objects stores, data parallel HDFS and Hbase (Bigtable), supporting close linking of computing and data location, as well as NOSQL table structures such as Cassandra and commercial approaches such as Amazon SimpleDB and Azure Table. Integration of CDMI standards as well as object repository standards.
- Economic models for an ecosystem with multiple cloud provides and services as well as publicly funded CI.
- Research on Cloud software stacks. There includes research at all levels of the software stack with two rather different emphasis areas: (1) research on systems that provide basic virtual machine provisioning, deployment and management, such as for example, Eucalyptus, Nimbus, OpenStack and OpenNebula with virtual networking as a distinct activity; and (2) integration of capabilities to provide rich Platform-as-a-Service as offered by major commercial systems – concepts such as appliances provide novel ways of delivering these capabilities.
- Clouds tend to achieve scalability by trading off overall reliability and robustness. Research is needed on both, how to expose faults to users as well as services to build fault tolerant applications. Most research in HPC tends to assume the absence of faults; however Clouds highlight a different philosophy with resilient applications running on faulty systems.
- Energy efficiency and Green IT is naturally synergistic with Clouds and related research includes examining the impact of Cloud features on power use, including the cost of powering idle machines supporting elastic clouds as well as a application aware approaches to power management.
- Performance modeling, engineering of Clouds and Cloud applications. In order to diagnose and tune performance, we need instrumentation at multiple levels, in both data and control planes: at the edge, in general infrastructure, in local clusters, in the control plane, and all the way down into the *aaS services. In addition, we also need benchmarks that are lightweight yet conclusive for ongoing tracking of cloud offerings that would for example, allow on-the-fly comparisons.
- End-to-end networking and data transport infrastructure and network virtualization.

4. Security Policies and Mechanisms

Clouds tend to emphasize the need for quality security mechanisms due to the sharing of storage and computing. One research area investigates hybrid architectures with algorithms broken into two; a low cost but non privacy preserving part running on an intrinsically secure private clouds, and a time consuming but privacy preserving part executing on a public cloud. Genomic data (human) and other health records are demanding here. The concept of differential privacy and health data anonymization is an active research topic. In addition to basic security for computing and storage there is research on privacy preserving search with the elegant but time consuming concept of Homomorphic Encryption, which allows encrypted data to be searched by encrypted queries. Some key research areas were highlighted at the recent NSF Workshop on Security for Cloud Computing³ include:

Adversary models for cloud computing:

- Identification of new security threats.
- Identification of possible sources of attacks – due to the different roles (user, providers, etc.) in cloud computing.
- Understanding how attacks in cloud computing are organized.
- Pricing and categorizing the security level.
- Relying on delegated technology and underlying cloud technologies.

Delegation and authorization in cloud computing:

- Delegation and authorization – such as attribute-based encryption for access control, secure comparison for complex policy enforcement, and encryption delegation for fine-grained temporal context.
- Mobile device access.
- Computation over encrypted data – when utilizing homomorphic encryption and homomorphic signatures.
- End-to-end cloud life cycle: restricted delegation, secure service composition, multiple credential types, and fine-grained access control.
- Restriction delegation (authorization based on capabilities not based on typical identification).

End-to-end security in cloud computing:

- Verification of work on clouds on behalf of clients.
- Kind of checks required at client and cloud sides.
- Trust between client and cloud.
- Security as a service within a cloud.
- Policy-based security applied to the end-to-end problem.
- Dealing with end-to-end privacy.
- Data ownership in the cloud.
- Root of trust - hardware, data, software (hypervisor), infrastructure (network) level.

New problems in security for cloud computing:

³ NSF Workshop on Security for Cloud Computing
<http://illinois.edu/blog/view/695/66281?count=1&ACTION=DIALOG>

- Considerations of legal service level agreements (SLAs) and stronger privacy policies for content providers who collect large amount of personal data.
- Attestation of mechanisms in clouds – e.g., trusted launch of VM and VM migration.
- Attestation of actions and proof-attestation of provider security mechanisms.
- Execution of algorithms on encrypted data.
- Cloud forensics – e.g., secure and correlate temporal and spatial evidence, use of log-based event for reconstruction.
- Reactive stability challenges, cross-layer robustness, pervasive virtualization, secure migration of data, storage, dependencies between services, placement, and management vulnerabilities.

In addition, challenges related to seamlessly using CI + Clouds as a single infrastructure include issues such as single sign-on, federated identify management (e.g., inCommons, cilogin, SCIM, etc.). Other issues include validation and security contextualization of virtual machines, establishing trust between virtual machines, tracking the provenance of virtual machines, etc.

5. Data Management in the Cloud

Cloud storage options. Cloud providers currently offer a large variety of storage options with varying service levels in terms of availability, access times, persistence, price, ease of use, etc. -- from ephemeral storage on virtual machines to storage clouds. As with other types of resource in the cloud it is hard to compare them and easily leverage the best solution for a given need.

Cloud connectivity. There is currently insufficient networking support for existing clouds, especially commercial offerings: moving data in takes a lot of time and is not happening fast enough. Fast networks, such as the CENIC network help alleviate the problem to a certain extent and for specific communities; the performance of those networks and their interaction with cloud storage offerings however has been little studied to date.

Data Locality. The problem of combining data and computation in one place is in many way similar to the same problem in grids but more complex due to a larger variety of storage options and connectivity constraints.

6. Deployment/Transition to Practice

Standards: There are many important standard activities, from those specifying the basic virtual machine structure to higher-level standards defining the PaaS environment, for example, queue and table structures. Although there is some support for these standards – such as OCCI (from OGF) in OpenNebula and OpenStack, and work at DMTF (dmtf.org/standards/cloud) – this area is still under development. NIST and IEEE are playing leadership roles.

Procurement and Account Management Tooling: One of the major obstacles today is the state of account management and procurement options for obtaining and using cloud services. One of

the issues is that resource management and usage in scientific community is oriented around a “wholesale” model, where scientists invest in and operate large computing centers funded by large grants, rather than a “retail” model where scientists lease resources, potentially for a short time and paying small amounts of money, effectively outsourcing their operation. A related problem deals with lack of established vehicles for the management of finely differentiated offering (e.g., acknowledging the difference between on-demand and best effort availability). Substantial progress needs to be made in order enable procurement of services through capped purchase orders, or subcontracts; administration of sub-accounts and delegation of resources and authorities, reporting and alarming/limiting (e.g., prohibiting specific runs from using more than a certain dollar amount at a time), etc.

Lack of Appropriate Middleware and Modes of Usage. Some examples were described in the preceding sections; they include lack of appropriate security tools and practices (e.g., how is intellectual property handled in the clouds?), lack of basic infrastructure ensuring cloud bursting and reliability, lack of comparison/monitoring tools and others.

Documentation: Documentation on how to use and the details of the system specific to scientific community internals are both inadequate. Development work is needed in partnership with the cloud providers, to improve the general documentation and particularly “HOW-TO” use cases for research usage.

Technical Training: The development and/or modification of codes adapted to the cloud environment require a set of skills currently in very short supply. These skills need to be taught, particularly to new practitioners in research, but also to mid-career practitioners in order to extract full benefit of this new technology.

Prototyping Support and Testbeds: Community testbeds that enable end-to-end research and experimentation, facilitate sharing of results and software, and well as support training and education, can be critical drivers.

Ease of Use: Much of the complexity associated with cloud computing could be simplified by the development of tools targeting scientific projects, adapting existing infrastructure to this new technologies. Most requested methods include infrastructure/instruction for adaptation of existing key management methods, better automation, and reduction in deployment times.

7. Research Timelines

Key **short term** activities include: (1) development of pilot projects to identify use cases and best practices, as well as to determine requirements; (2) definition of standards and develop services to enable interoperability between CI (i.e., XSEDE, OSG, FutureGrid) and Cloud services (e.g., EC2, Azure, etc.) as well as their integration; (3) creation of community research and experimentation testbeds as well as establishment community benchmarks and common metrics; (4) translation of research innovation into software frameworks tools that can be used by applications; and (5) creation of community forums for exchange of ideas and artifacts. **Longer-term** activities include addressing the fundamental research challenges discussed above and

incorporating resulting research innovations into sustainable software system that can be deployed and used by the community. Establishing appropriate research program and community forums to support these research activities will be important.